# Assignment No. 4

## 4.1 Title

Consider a suitable text dataset. Remove stop words, apply stemming and feature selection techniques to represent documents as vectors. Classify documents and evaluate precision, recall.

## 4.2 Problem Definition:

Remove stop words

## 4.3 Prerequisite:

Basic Concepts of ETL

## 4.4 Software Requirements:

Rapid Miner

## 4.5 Hardware Requirement:

PIV, 2GB RAM, 500 GB HDD, Lenovo A13-4089Model.

## 4.6 Learning Objectives:

We are going to learn how to tokenize and filter a document into its different words and then do words count for each word in a text document

## 4.7 Outcomes:

You are able to see a word list containing all the different words in your document and their occurrence count next to it in the "Total Occurrences" column
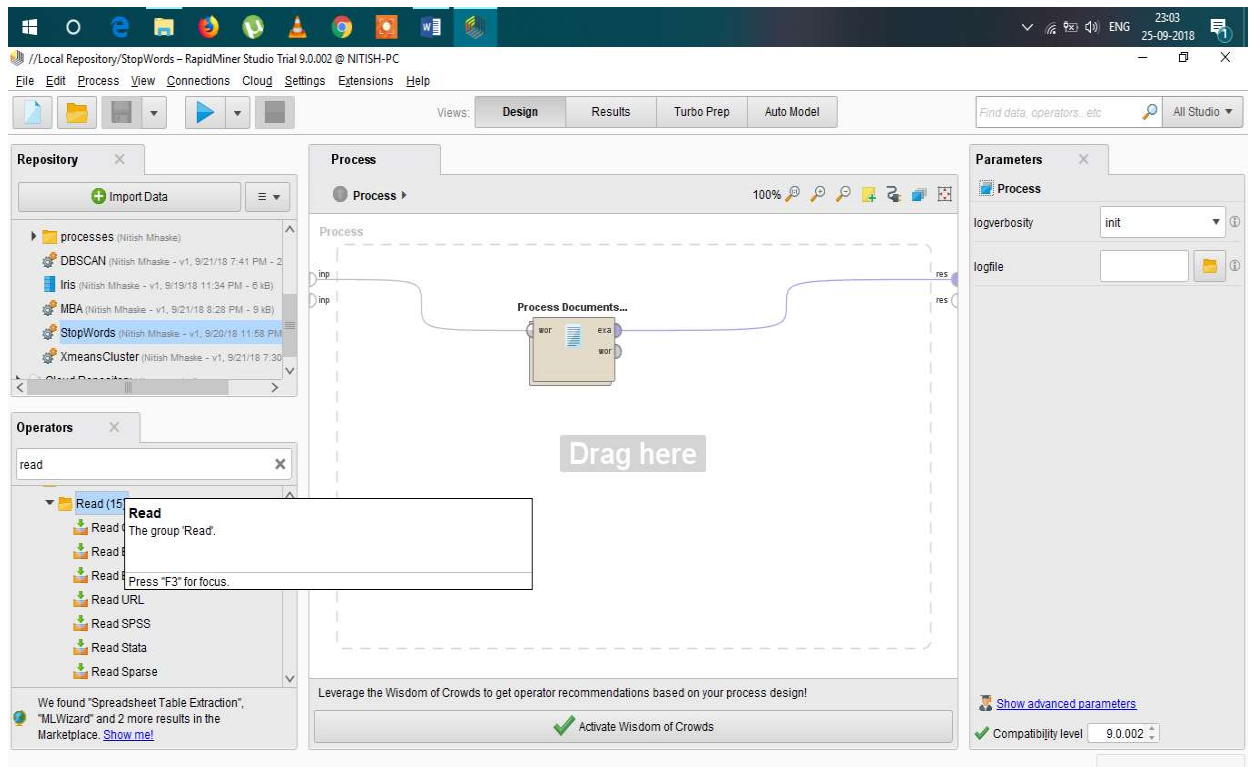
## 4.8 Theory Concepts:

### 4.8.1 Text Processing Tutorial with RapidMiner

In this, we are going to learn how to tokenize and filter a document into its different words and then do words count for each word in a text document.
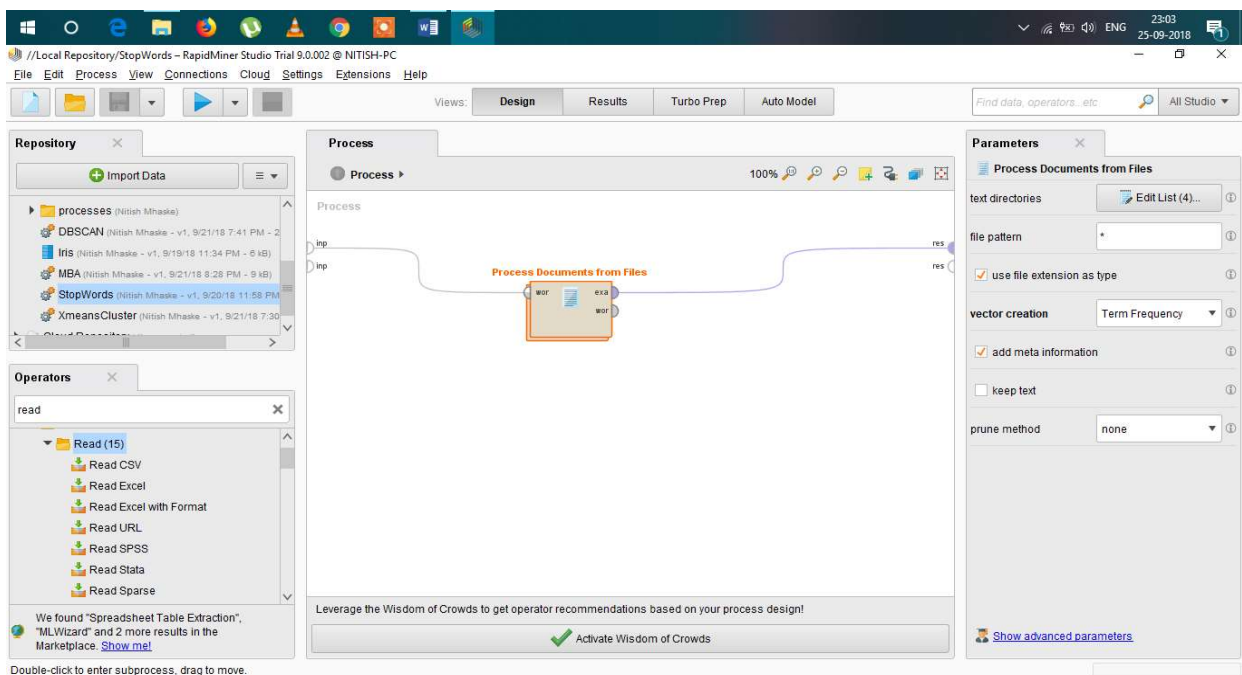
Steps:

1. Open RapidMiner and click "New Process". On the left hand pane of your screen, there should be a tab that says "Operators"- this is where you can search and find all of the operators for RapidMiner and its extensions. By searching the Operators tab for "read".
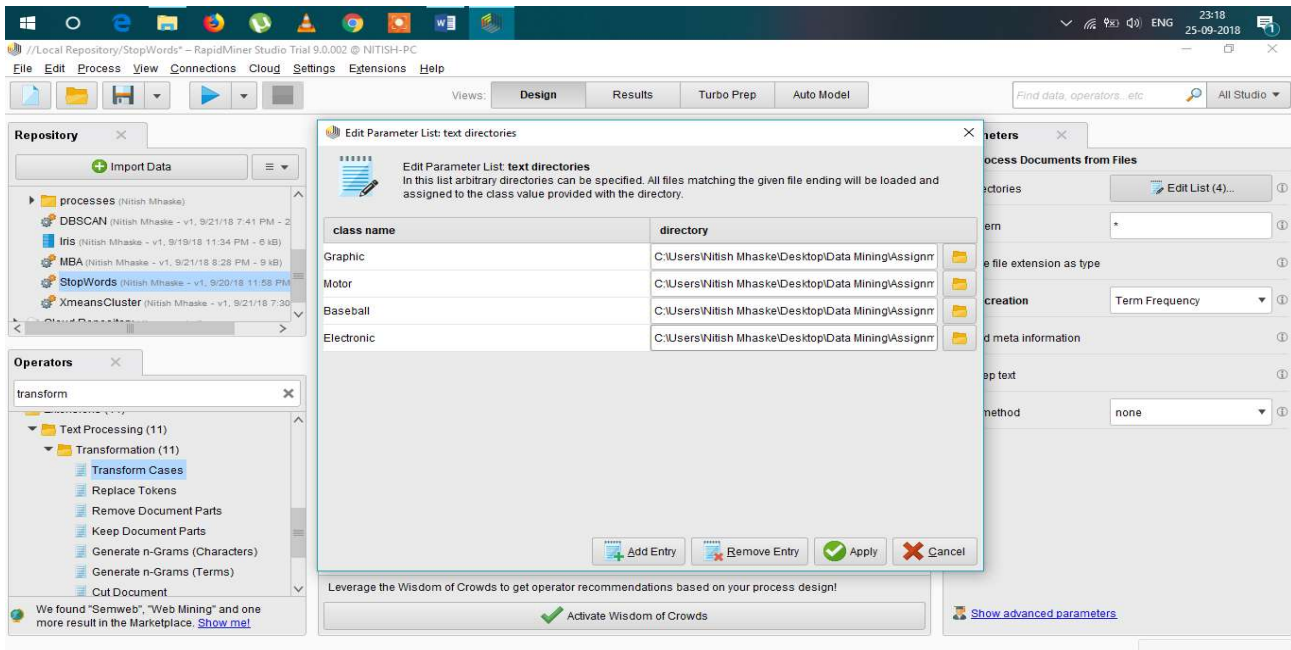
**Figure 4.1:** Searching the Operators tab for "read"

2. There are multiple read operators depending on which file you have, and most of them work the same way. If you scroll down, there is a "Read Documents" operator. Select this operator and enter it into your Main Process window by dragging it. When you select the Read Documents operator in the Main Process window, you should see a file up loader in the right-hand pane.
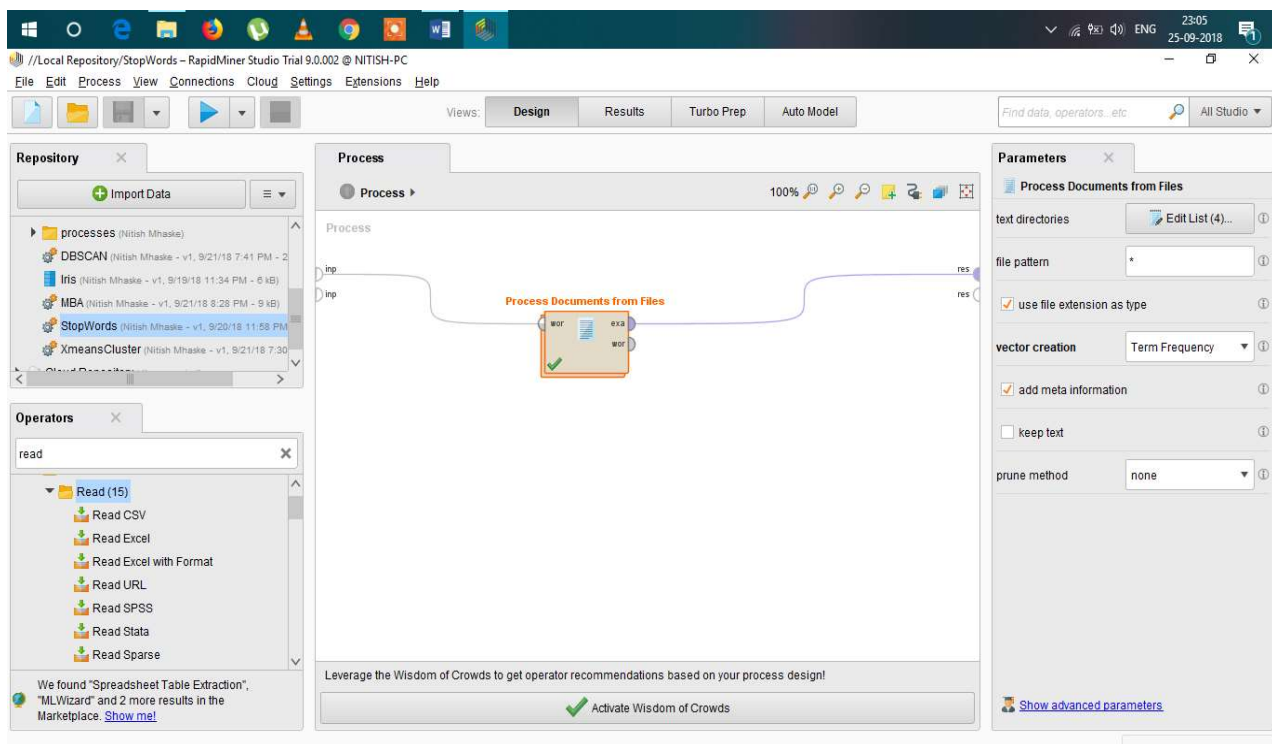


**Figure 4.2:** Drag and Drop "Process Documents" operator

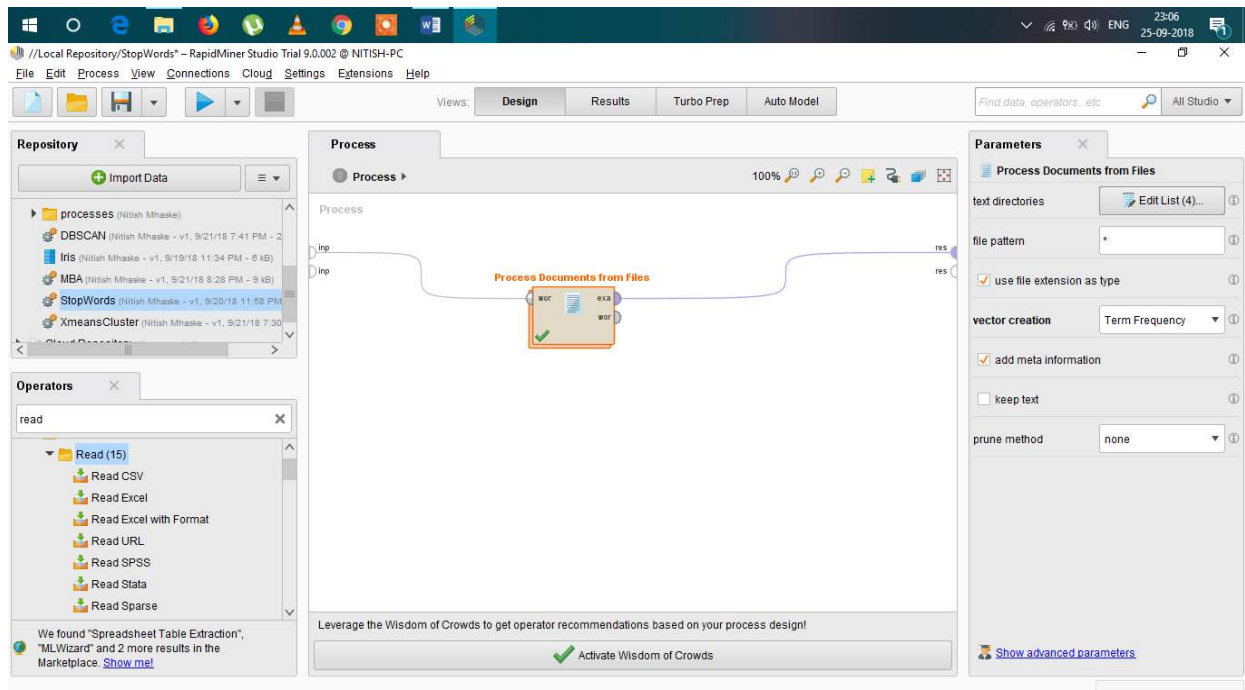3.  Select the text file you want to use.



**Figure 4.3:** Select the text file you want to use

4.  After you have chosen your file, make sure that the output port on the Read Documents operator is connected to the "res" node in your Main Process. Click the "play" button to check that your file has been received correctly. Switch to the results perspective by clicking the icon that looks like a display chart above the "Process" tab at the top of the Main Process pane. Click the "Document (Read Document)" tab. Your output text should look something like this depending on the file you have chosen to process:
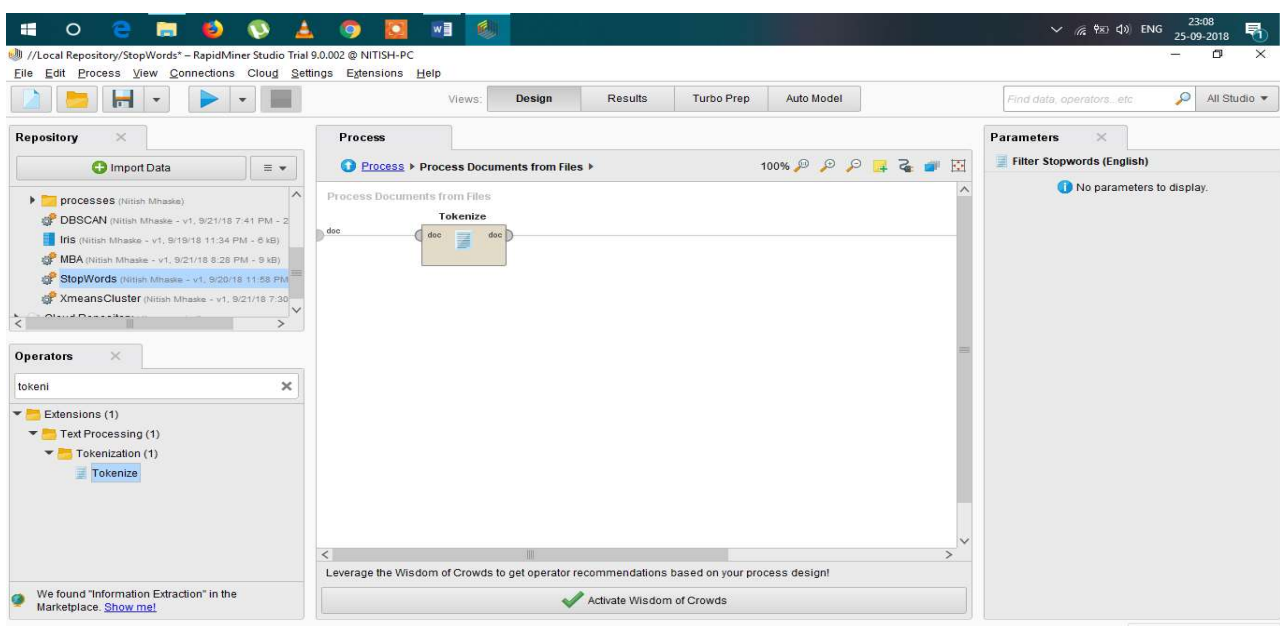


**Figure 4.4:** Run the Process

5. Now we will move on to processing the document to get a list of its different words and their individual count. Search the Operators list for "Process Documents". Drag this operator the same way as you did for the "Read Documents" operator into the main panel.
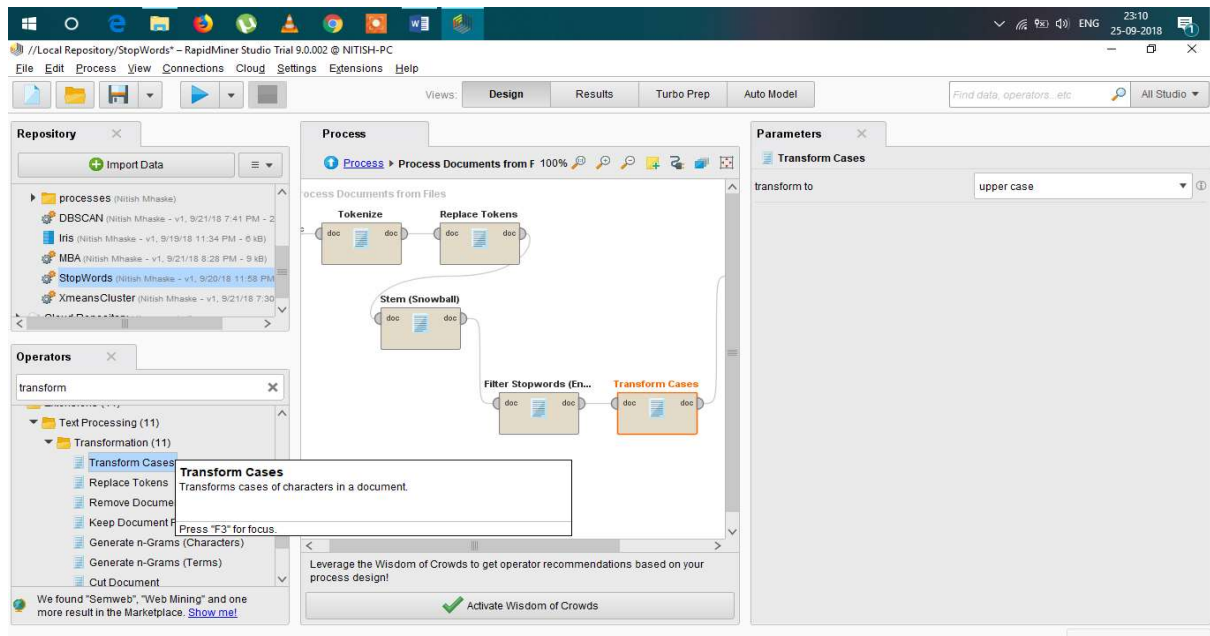


**Figure 4.5:** Search the Operators list for "Process Documents"

6. Double click the Process Documents operator to get inside the operator. This is where we will link operators together to take the entire text document and split it down into its word components. This consists of several operators that can be chosen by going into the Operator pane and looking at the Text Processing folder. You should see several more folders such as "Tokenization", "Extraction", "Filtering", "Stemming", "Transformation", and "Utility".
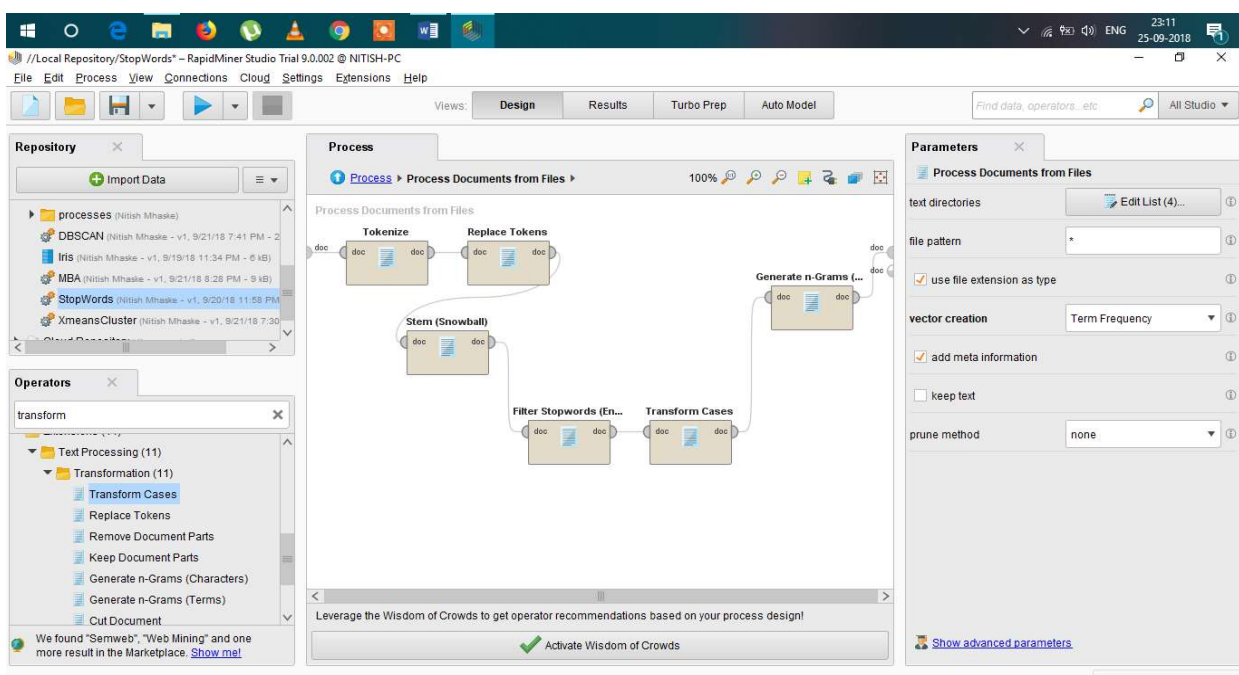


**Figure 4.6** Search for the "Tokenize" operator

7. Connect the "doc" node of the process to the "doc" input node of the operator if it has not automatically connected already. Now we are ready to filter the bag of words. In "Filtering" folder under the "Text Processing" operator folder, you can see the various filtering methods that you can apply to your process.



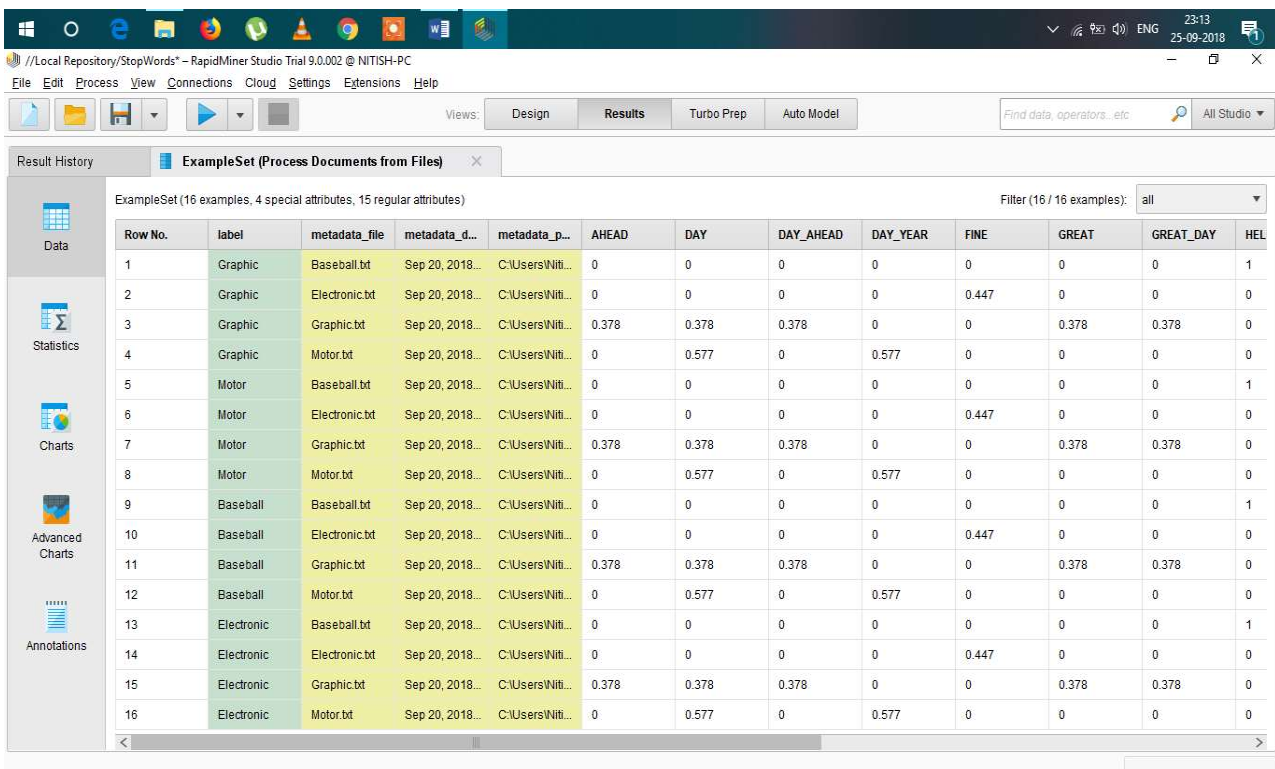**Figure 4.7** Select the operator "Transform Cases" and drag it into the process.

8. After I filtered the bag of words by stopwords and length, I want to transform all of my words to lowercase since the same word would be counted differently if it was in uppercase vs. lowercase. Select the operator "Transform Cases" and drag it into the process.



**Figure 4.8** Checks all nodes connections and clicks the "Play" button to run process

9.  Now that I have the sufficient operators in my process for this example, I check all of my node connections and click the "Play" button to run my process. If all goes well, your output should look like this in the results view:



**Figure 4.9** Output should look like this in the results view

## 4.9 Conclusion

  We are now able to see a word list containing all the different words in your document and their occurrence count next to it in the "Total Occurrences" column. If you do not get this output, make sure that all of your nodes are connected correctly and also to the right *type*. Some errors are because your output at one node does not match the type expected at the input of the next node of an operator

**References:-** https://docs.rapidminer.com/downloads/RapidMiner-v6-user-manual.pdf